

Smoothing, Clustering, and Benchmarking for Small Area Estimation

Rebecca C. Steorts, Carnegie Mellon University

October 28, 2014

Abstract

We develop constrained Bayesian estimation methods for small area problems: those requiring smoothness with respect to similarity across areas, such as geographic proximity or clustering by covariates; and benchmarking constraints, requiring (weighted) means of estimates to agree across levels of aggregation. We develop methods for constrained estimation decision-theoretically and discuss their geometric interpretation. Our constrained estimators are the solutions to tractable optimization problems and have closed-form solutions. Mean squared errors of the constrained estimators are calculated via bootstrapping. Our techniques are free of distributional assumptions and apply whether the estimator is linear or non-linear, univariate or multivariate. We illustrate our methods using data from the U.S. Census’s Small Area Income and Poverty Estimates program.

1 Introduction

Small area estimation (SAE) deals with estimating many parameters, each associated with an “area”—a geographic domain, a demographic group, an experimental condition, etc. Areas are “small” since there is little or no information about any one area. Estimates of a parameter based only on observations from the associated area, called direct estimates, are imprecise. To increase precision, one tries to “borrow strength” from related areas, and hierarchical and empirical Bayesian models are one way to do so. Since the pioneering work of [Battese et al. \(1988\)](#); [Fay and Herriot \(1979\)](#), such models have dominated SAE, with many successful applications in official statistics, sociology, epidemiology, political science, business, etc. ([Rao 2003](#)). Recently, SAE has been applied in other fields, such as

neuroscience, where it has been shown to do as well as common approaches such as smoothed ridge regression and elastic net (Wehbe et al. 2014).

We extend these classical approaches in two directions, both of which have been the subject of recent interest in the SAE literature. One direction is to directly take account of information about the proximity of areas in space or time. In many applications, it is reasonable to expect that the parameters will be smooth, so that nearby areas will have similar parameters, but this is not altogether standard within SAE (Rao 2003). Incorporating spatial or temporal dependence directly into Bayesian models leads to statistical and computational difficulties, yet it seems misguided to discard such information. The other direction is “benchmarking,” the imposition of consistency constraints on (weighted) averages of the parameter estimates. A simple form of benchmarking is when the average of the parameter estimates must match a known global average. When there are multiple levels of aggregation for the estimates, there can be issues of internal consistency as well.

We provide a unified approach to smoothing and benchmarking by regarding them both as *constraints* on Bayes estimates. Benchmarking corresponds to equality constraints on global averages and variances. Similarly, smoothing corresponds to an inequality constraint on the “roughness” of estimates (how much the parameter estimates of nearby areas differ). The motivation of this smoothing is based upon manifold learning and frequentist non-parametrics, where loss functions are augmented by a penalty. Such a penalty term is in the spirit of ridge regression, where a transformation of the parameters is performed and additional shrinkage is carried out. Our penalty corresponds to how much estimates at nearby points in the domain should tend to differ.

Decision-theoretically, we obtain smoothed, benchmarked estimates by minimizing the Bayes risk subject to these constraints, extending the approaches of Datta et al. (2011); Ghosh and Steorts (2013) (themselves in the spirit of Louis (1984) and Ghosh (1992)). Geometrically, the constrained Bayes estimates are found by projecting the unconstrained estimates into the feasible set. If the constraints are linear, then the the resulting optimization can be solved in closed form, requiring nothing more than basic matrix operations on the unconstrained Bayes estimates. When we use equality constraints that are quadratic in nature, the problem cannot be solved in closed form, and the optimization is in fact non-convex.

Previous efforts at smoothing in SAE problems have smoothed either the raw data or direct estimates. In contrast, we smooth estimates based on models which do not themselves include spatial structure. Computationally,

this is much easier than expanding the models. Our optimization problems can be solved in closed form and retain the advantages of model-based estimation. This approach to smoothing also combines naturally with the imposition of benchmarking constraints.

Another strong advantage of our decision-theoretic and geometric approach is its extreme generality. We require no distributional assumptions on the data or on the unconstrained Bayes estimator. Our results apply whether the unconstrained estimator is linear or non-linear, whether the parameters being estimated are univariate or multivariate, and whether there is a single level of aggregation (“area”-level models) or multiple (“unit”-level models). The relevant notion of proximity between areas may be spatial, temporal, or more abstract. It can even include clustering on covariates not directly included in the Bayesian model.

§2 describes our approach to smoothing in small area estimation. This is extended to add benchmarking constraints in §3. In both sections, area- and unit-level results are derived using a single framework. §4 discusses uncertainty quantification based on a bootstrap approach. §5 presents an application to estimating the number of children living in poverty in each state, using data from the U.S. Census Bureau’s Small Area Income and Poverty Estimates Program (SAIPE).

1.1 Notation and Terminology

We assume m **areas**, and for each area i , we estimate an associated scalar quantity θ_i , collectively $\boldsymbol{\theta}$. “Areas” are often spatial regions, where they might be different demographic groups or experimental conditions. Allowing θ_i to be vectors rather than scalars is straightforward (see the remark at the end of §2). Each area has a vector of covariates \boldsymbol{x}_i , which may include spatial or temporal coordinates, when applicable. Conditioning a Bayesian model on an observed response \boldsymbol{y} and covariates \boldsymbol{x} leads to a Bayes estimate $\hat{\theta}_i^B$ for each area. (Note that $\hat{\theta}_i^B$ is obtained by conditioning on *all* the observations and covariates, not just those of area i .)

The loss function is weighted squared error, where the weight for area i is $\phi_i > 0$, and the total loss from the action (estimate) $\boldsymbol{\delta}$ is $\sum_i \phi_i (\theta_i - \delta_i)^2$. In many SAE applications, ϕ_i reflect variations in measurement precision and can be obtained from the survey design (Pfeffermann 2013; Rao 2003). We assume they are known (however, in practice they usually must be estimated). Define $\Phi = \text{Diag}(\phi_i)$, which is positive definite by construction.

1.2 Related Work

Pfeffermann (2013) provided a comprehensive review of the SAE and benchmarking literature. Our work is twofold: smoothing SAE, and its combination with benchmarking. Our SAE approach is decision-theoretic with the addition of a *smoothness penalty* in the loss function. Our approach to benchmarking with smoothing generalizes the benchmarking work of Datta et al. (2011); Ghosh and Steorts (2013).

It is thought that spatial correlations *may* help SAE models, leading to approaches such as correlated sampling error, spatial dependence of small area effects, time series components, etc., reviewed in Ghosh and Rao (1994). More recently, spatially-correlated random effects have been incorporated into empirical Bayesian models (Pratesi and Salvati 2008) and into hierarchical Bayesian models (Souza et al. 2009). These approaches have all been highly application-specific and hard to integrate with benchmarking, and they greatly increase the computational cost of obtaining estimates. Our goal is to overcome these limits by taking a radically different approach.

Thus, we employ ideas about smoothing on graphs and manifolds from frequentist non-parametrics and machine learning. In particular, we take advantage of “Laplacian” regularization ideas (Belkin et al. 2006; Corona et al. 2008; Lee and Wasserman 2010), where the loss function is augmented by a penalty term which reflects how much estimates at nearby points in the domain should tend to differ. Such regularization is designed to ensure that estimates vary smoothly with respect to the intrinsic geometry of some underlying graph or manifold. (Smoothness on a domain is represented mathematically by the domain’s Laplacian operator, which is the generator for diffusion processes.) This generalizes the roughness or curvature penalties from spline smoothing (Wahba 1990) to domains more geometrically complicated than \mathbb{R}^m . We are unaware of any previous application of Laplacian regularization to SAE problems, though spline smoothing is often used in spatial statistics, including such classic SAE tasks as estimating disease rates (Kafadar 1996).

2 Smoothing for Small Area Estimation

In this section, we develop estimators that minimize posterior risk while still imposing smoothness on the estimate. The kind of smoothing we impose derives from the literature on Laplacian regularization and semi-supervised learning (Belkin et al. 2006; Corona et al. 2008; Lee and Wasserman 2010). The estimators we derive do not depend on the distributional assumptions

of the Bayesian models and are equally applicable to spatial smoothing or more abstract clustering. We do assume that the smoothing or clustering is done separately from the estimation for each area or domain, and we also take weighted squared error as the loss function. In §3, we extend our approach to include benchmarking.

2.1 General Result

We begin by introducing the symmetric matrix Q , with elements $q_{ii'} \geq 0$, to gauge how important it is that the estimate of θ_i be close to the estimate of $\theta_{i'}$. It may often be the case that $q_{ii'} = q(\mathbf{x}_i, \mathbf{x}_{i'})$, i.e., the degree of smoothing of δ_i and $\delta_{i'}$ is a function of the covariates \mathbf{x}_i and $\mathbf{x}_{i'}$. Note also that the $q_{ii'}$ may be discrete-valued, corresponding to clustering of areas, or continuous-valued, corresponding to a metric space of areas.

A natural measure of the smoothness of $\boldsymbol{\delta}$ is the Q -weighted sum of squared differences between elements, $\sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'}$. Hence, we add a penalty term $\gamma \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'}$ to our objective function, with the penalty factor $\gamma \geq 0$ chosen to specify the overall importance of smoothness. (We address the choice of Q below and of γ in §3.3.)

Therefore, we seek to minimize the posterior risk of the loss function

$$L(\boldsymbol{\theta}, \boldsymbol{\delta}) = \sum_i \phi_i(\theta_i - \delta_i)^2 + \gamma \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'}. \quad (1)$$

Minimizing the posterior expectation of (1) is equivalent to minimizing

$$\sum_i \phi_i E[(\theta_i - \delta_i)^2 | \mathbf{y}] + \gamma \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'}. \quad (2)$$

Define Ω to be a matrix such that $\sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'} = \boldsymbol{\delta}^T \Omega \boldsymbol{\delta}$. (See Lemma 1 for details.) Then minimizing (2) is equivalent to minimizing

$$(\boldsymbol{\delta} - \hat{\boldsymbol{\theta}}^B)^T \Phi (\boldsymbol{\delta} - \hat{\boldsymbol{\theta}}^B) + \gamma \boldsymbol{\delta}^T \Omega \boldsymbol{\delta}. \quad (3)$$

(We refer to (Datta et al. 2011; Ghosh and Steorts 2013) for details on this equivalence.) Then we have the following result.

Theorem 1. *The smoothed Bayes estimator is*

$$\tilde{\boldsymbol{\theta}}^S = (I_m + \gamma \Phi^{-1} \Omega)^{-1} \hat{\boldsymbol{\theta}}^B.$$

Proof. Differentiating (3) and setting the gradient to zero at $\tilde{\boldsymbol{\theta}}^S$ yields $\Phi(\tilde{\boldsymbol{\theta}}^S - \hat{\boldsymbol{\theta}}^B) + \gamma\Omega\tilde{\boldsymbol{\theta}}^S = \mathbf{0}$. Then

$$(\Phi + \gamma\Omega)\tilde{\boldsymbol{\theta}}^S = \Phi\hat{\boldsymbol{\theta}}^B \implies \tilde{\boldsymbol{\theta}}^S = (I_m + \gamma\Phi^{-1}\Omega)^{-1}\hat{\boldsymbol{\theta}}^B.$$

Since (3) is a positive-definite quadratic form in $\boldsymbol{\delta}$, the solution is unique. \square

See §A.1 for an extension to unit-level models.

Remark. The parameter to be estimated for each area may be multivariate. For instance, we might seek both a poverty rate and a median income for each area. For simplicity, we assume that the parameter dimension p is the same for each of the m areas. Then Theorem 1 can be applied with $\boldsymbol{\theta} = (\theta_{11}, \dots, \theta_{m1}, \dots, \theta_{12}, \dots, \theta_{m2}, \dots, \theta_{1p}, \dots, \theta_{mp})$. The matrix Φ remains the diagonal matrix of the ϕ_{ij} , in the same order as $\boldsymbol{\theta}$. However, Ω is now a block-diagonal matrix, where each $m \times m$ block contains a copy of the appropriate matrix for the corresponding univariate problem. This ensures that the same smoothness constraint is imposed on each component of the parameter vectors, but different components are not smoothed together.

3 Benchmarking and Smoothing

We now turn to situations where our estimates should not just be smooth, minimizing (3), but also obey benchmarking constraints. As the benchmarking constraints are relaxed, we should recover the results of §2. Our approach to finding benchmarked Bayes estimators extends that of Datta et al. (2011); Ghosh and Steorts (2013). We employ the following definition.

Definition 1 (Benchmarking constraints, benchmarked Bayes estimator). *Benchmarking constraints* are equality constraints on the weighted means or weighted variances of subsets (possibly all) of the estimates. The *benchmarking Bayes estimator* is the minimizer of the posterior risk subject to the benchmarking constraints.

The levels to which we benchmark, i.e., the values of the equality constraints, are assumed to be given *externally* from some other data source. (For internal benchmarking, see Bell et al. (2013).) Our methods address linear, weighted mean constraints, as in Datta et al. (2011); Ghosh and Steorts (2013).

3.1 Linear Benchmarking Constraints

We first consider benchmarking constraints which are linear in the estimate $\boldsymbol{\delta}$, such as means or totals. The general problem is now to minimize the posterior risk in (3) subject to the constraints

$$M\boldsymbol{\delta} = \mathbf{t}, \quad (4)$$

where \mathbf{t} is a given k -dimensional vector and M is a $k \times m$ matrix. As before, this is equivalent to introducing a k -dimensional vector of Lagrange multipliers $\boldsymbol{\lambda}$ and minimizing

$$(\boldsymbol{\delta} - \hat{\boldsymbol{\theta}}^B)^T \Phi (\boldsymbol{\delta} - \hat{\boldsymbol{\theta}}^B) + \gamma \boldsymbol{\delta}^T \Omega \boldsymbol{\delta} - 2\boldsymbol{\lambda}^T (M\boldsymbol{\delta} - \mathbf{t}).$$

Theorem 2. Suppose (4) has solutions. Then the constrained Bayes estimator under (4) is

$$\tilde{\boldsymbol{\theta}}^{BM} = \Sigma^{-1} \left[\Phi \hat{\boldsymbol{\theta}}^B + M^T (M \Sigma^{-1} M^T)^{-1} (\mathbf{t} - M \Sigma^{-1} \Phi \hat{\boldsymbol{\theta}}^B) \right],$$

where $\Sigma = \Phi + \gamma \Omega$.

Remark. Note that the Theorem 2 estimator $\tilde{\boldsymbol{\theta}}^{BM}$ can be expressed in terms of the Theorem 1 estimator $\tilde{\boldsymbol{\theta}}^S$ as

$$\tilde{\boldsymbol{\theta}}^{BM} = \tilde{\boldsymbol{\theta}}^S + \Sigma^{-1} M^T (M \Sigma^{-1} M^T)^{-1} (\mathbf{t} - M \tilde{\boldsymbol{\theta}}^S).$$

Thus, it can be seen that the benchmarking essentially “adjusts” the estimator $\tilde{\boldsymbol{\theta}}^S$ based on the discrepancy between $M \tilde{\boldsymbol{\theta}}^S$ and the target \mathbf{t} .

Proof of Theorem 2. Differentiating with respect to $\boldsymbol{\delta}$ and setting the result equal to zero at $\tilde{\boldsymbol{\theta}}^{BM}$ yields

$$\begin{aligned} M^T \boldsymbol{\lambda} &= \Phi (\tilde{\boldsymbol{\theta}}^{BM} - \hat{\boldsymbol{\theta}}^B) + \gamma \Omega \tilde{\boldsymbol{\theta}}^{BM} \\ \implies \tilde{\boldsymbol{\theta}}^{BM} &= \Sigma^{-1} (\Phi \hat{\boldsymbol{\theta}}^B + M^T \boldsymbol{\lambda}). \end{aligned}$$

Then by the constraint,

$$\begin{aligned} \mathbf{t} &= M \Sigma^{-1} (\Phi \hat{\boldsymbol{\theta}}^B + M^T \boldsymbol{\lambda}) \\ &= M \Sigma^{-1} \Phi \hat{\boldsymbol{\theta}}^B + M \Sigma^{-1} M^T \boldsymbol{\lambda}, \end{aligned} \quad (5)$$

so $\boldsymbol{\lambda} = (M \Sigma^{-1} M^T)^{-1} (\mathbf{t} - M \Sigma^{-1} \Phi \hat{\boldsymbol{\theta}}^B)$. The result follows immediately. \square

Often there is only one linear constraint of the form $\sum_i w_i \delta_i = t$, or equivalently $\mathbf{w}^T \boldsymbol{\delta} = t$, for some nonnegative weights w_i and some $t \in \mathbb{R}$. This is simply a special case of Theorem 2 with $k = 1$ and $M = \mathbf{w}^T$, in which case the result simplifies to

$$\tilde{\boldsymbol{\theta}}^{BM} = \tilde{\boldsymbol{\theta}}^S + (t - \mathbf{w}^T \tilde{\boldsymbol{\theta}}^S) (\mathbf{w}^T \Sigma^{-1} \mathbf{w})^{-1} \Sigma^{-1} \mathbf{w}.$$

Also see §A.2.1 for an extension to unit-level models.

3.2 Weighted Variability Constraints

Now suppose that we impose one additional constraint on the weighted variability of the form

$$(\boldsymbol{\delta} - \boldsymbol{\tau})^T W (\boldsymbol{\delta} - \boldsymbol{\tau}) = h \quad (6)$$

where $\boldsymbol{\tau}$ is an m -dimensional vector, W is a symmetric $m \times m$ matrix, and h is a non-negative constant. In the important special case of a single weighted-mean constraint of the form $\boldsymbol{w}^T \boldsymbol{\delta} = t$, the vector $\boldsymbol{\tau}$ and the matrix W in (6) are typically taken as $\boldsymbol{\tau} = t \mathbf{1}_m$ and $W = \text{Diag}(w_i)$, where $\mathbf{1}_m$ denotes the ones vector of length m .

The imposition of such a constraint immediately renders the associated optimization problem considerably more challenging. Specifically, closed-form solutions for the minimizer can no longer be obtained in general. Moreover, notice that the set of $\boldsymbol{\delta} \in \mathbb{R}^m$ that satisfy both (4) and (6) is no longer convex. Thus, even a numerical solution may be difficult to obtain. The question of how to incorporate variability constraints while maintaining tractability of the model is a potential direction of future research and is beyond the scope of this paper.

3.2.1 Geometric Interpretation

Our formulation of benchmarking and smoothing as constrained optimization problems has a very simple geometric interpretation. It is well known that the Bayes estimate is the minimizer of the conditional expectation of the mean squared error (MSE). Since the minimization is taken over *all* possible values of $\boldsymbol{\theta}$, the Bayes estimate will not respect any constraints we might wish to impose (except by chance) or unless these constraints are included in the specification of the prior. We instead seek to minimize the MSE within the feasible set of the constraints. We find the point in the feasible set which is as close (in the sense of expected weighted squared error) to the Bayes estimate as possible. That is, we *project* the Bayes estimate into the feasible set.

The geometry of the feasible set is itself slightly complicated, because of the constraints imposed. Note that the smoothness penalty in the loss function may be reformulated as a smoothness constraint of the form $\boldsymbol{\delta}^T \Omega \boldsymbol{\delta} \leq s$ for some $s > 0$. This constraint defines an ellipsoid centered at the origin. Constraints on weighted means define linear sub-spaces, e.g., planes, depending on the number of constraints and the number of variables. Finally, constraints of weighted variabilities define the surfaces of cones. The

constrained Bayes estimator is the projection of the unconstrained Bayes estimator onto the intersection of the ellipsoid, the linear sub-space, and the cones.

3.3 Choice of Smoothing Penalties

The choice of γ ¹ is assumed fixed *a priori*. But knowing γ is equivalent to knowing how smooth the estimate *ought* to be, and this knowledge is lacking in most applications. In such situations, we suggest obtaining γ by leave-one-out cross-validation (Corona et al. 2008; Stone 1974; Wahba 1990).

For each value of γ and each area i , define $\boldsymbol{\delta}^{(-i)}(\gamma)$ as the solution of the corresponding optimization problem with the loss-function term for i dropped². The smoothness penalty and any applicable benchmarking constraints are however calculated over the *whole* of the vector $\boldsymbol{\delta}$, not just the non- i entries. (This ensures that $\boldsymbol{\delta}^{(-i)}(\gamma)$ does meet all the constraints, while still making a *prediction* about θ_i .) The cross-validation score of γ is

$$V(\gamma) = \frac{1}{m} \sum_{i=1}^m \left[\delta_i^{(-i)}(\gamma) - \hat{\theta}_i^B \right]^2 \phi_i,$$

where $\delta_i^{(-i)}(\gamma)$ denotes the i th component of $\boldsymbol{\delta}^{(-i)}(\gamma)$, and the minimizer of the cross-validation scores is $\hat{\gamma} = \operatorname{argmin}_{\gamma \geq 0} V(\gamma)$.

Direct evaluation of $V(\gamma)$ can be computationally costly. See Wahba (1990) for faster approximations, such as “generalized cross-validation.”

4 Evaluation Using a Residual Bootstrap

It is traditional in small area estimation to report approximations to the overall estimation error. (One of the main motivations of using small area methods is, after all, reducing the estimation error.) This is generally a challenging undertaking, since while methods like cross-validation can be used to evaluate *prediction* error in a way which is comparable across models, they do not work for *estimation* error. Thus, one needs to use more strictly model-based approaches, either analytic or based on the bootstrap.

Evaluating the MSE of our estimates is especially difficult, since we combine a model-based estimate with a non-parametric smoothing term.

¹In unit-level problems, γ_A and γ_U ; we will not note all the small modifications needed to pick two smoothing factors at once.

²Instead of the sum of squared errors $\sum_{i'=1}^m \phi_{i'}(\delta_{i'} - \theta_{i'})^2$, we use $\sum_{i' \neq i} \phi_{i'}(\delta_{i'} - \theta_{i'})^2$. This amounts to replacing Φ with a matrix whose i th row and column are both 0.

A straightforward model-based bootstrap would sample from the posterior distribution of (7) to generate a new set of true poverty rates θ^* and observations y^* , re-run the estimation on y^* , and see how close the resulting estimates δ^* came to θ^* . However, this presumes the correctness of the Fay-Harriot model in (7), which is precisely what we have chosen *not* to assume through our imposition of the benchmarking/smoothing constraints³. Note that such constraints do not fit naturally into the generative model.

We evade this dilemma by using a semi-parametric residual bootstrap, a common approach when the functional form of a regression is known fairly securely, but the distribution of fluctuations is not. We calculate the differences between y_i and our constrained Bayes estimates $\hat{\theta}_i^{BM}$, resample these residuals, scale them to account for heteroskedasticity, and add them back to $\hat{\theta}_i^{BM}$ to generate y_i^* . (See Appendix C.) The residual bootstrap assumes that smoothing is appropriate and that we have chosen the right Ω matrix.

5 Application to the SAIPE Dataset

We apply our constrained Bayes estimation procedure to data from the Small Area Income and Poverty Estimates (SAIPE) program of the U.S. Census Bureau. Our goal is to estimate, for 1998, the rate of poverty of children aged 5–17 years in each state and the District of Columbia. This is an area-level model, with states as the areas. The small area model from which we derive our initial Bayes estimates is described in §5.1. The primary benchmarking constraint is that the weighted mean of the state poverty-rate estimates must match the national poverty rate established by direct estimates. A secondary benchmarking constraint is the matching of the similarly-known national variance in poverty rates. This benchmarking had already been considered, for this data and model, in Datta et al. (2011). We add to this the constraint of smoothness across states, where our choice of Laplacian and of smoothing penalty is discussed in §5.2.

SAIPE estimates average household income and poverty rates from small areas in the U.S. over multiple years. It works by combining direct estimates of these quantities, from the Annual Social and Economic (ASEC) Supplement of the Current Population Survey (CPS) and the American Community Survey (ACS), with standard small area models, which use as predictors several variables drawn from administrative records. This presumes that areas with similar values of the predictor variables should have similar values of

³If we follow this procedure nonetheless, we always conclude that benchmarking and especially smoothing radically increase the MSE by introducing large biases.

the parameters of interest. See [Bell et al. \(2013\)](#); [Datta et al. \(2011\)](#) for a fuller account of the SAIPE program.

For illustrative purposes, and following [Datta et al. \(2011\)](#), we have focused on estimating the rate of poverty among children aged 5–17 in 1998. The public-use data here is at the state level, so states are areas. The predictor variables used are: a pseudo-estimate of the child poverty rate based on Internal Revenue Service (IRS) income-tax statements; the rate of households not filing taxes with the IRS⁴; the rate of food stamp⁵ use; and the residual term from a regression of the 1990 Census estimated child poverty rate. The last variable is supposed to help represent whether a state has an unusual level of poverty, given its other characteristics, which is presumably persistent over time.

5.1 Hierarchical Bayesian Model for SAIPE

As in [Datta et al. \(2011\)](#), our initial, unconstrained Bayes estimates for poverty rates are derived from the following hierarchical model of [Fay and Herriot \(1979\)](#):

$$\begin{aligned} y_i \mid \theta_i &\sim \mathcal{N}(\theta_i, D_i); \quad i = 1, \dots, m \\ \theta_i \mid \boldsymbol{\beta} &\sim \mathcal{N}(\mathbf{x}_i^T \boldsymbol{\beta}, \sigma_u^2) \\ \pi(\sigma_u^2, \boldsymbol{\beta}) &\propto 1. \end{aligned} \tag{7}$$

Here θ_i is the true poverty rate for state i , y_i is the direct survey estimate, D_i is the known sampling variance of the survey, \mathbf{x}_i are the predictors, $\boldsymbol{\beta}$ is the vector of regression coefficients, and σ_u^2 is the unknown modeling variance. The posterior means and variances of $(\boldsymbol{\beta}, \theta_i, \sigma_u^2)$ are estimated by Gibbs sampling.

5.2 Benchmarking and Smoothing Results

We consider three different possibilities for benchmarking and/or smoothing: (i) benchmarking the mean alone without smoothing, (ii) benchmarking both the mean and variability without smoothing, and (iii) benchmarking the mean alone with smoothing. Note that since there is no smoothing in (ii), solutions can indeed be found in closed form; see [Datta et al. \(2011\)](#) for details.

⁴In the U.S., households whose income falls below certain thresholds are not required to file federal taxes.

⁵A program providing direct assistance in buying food and other necessities for low-income households.

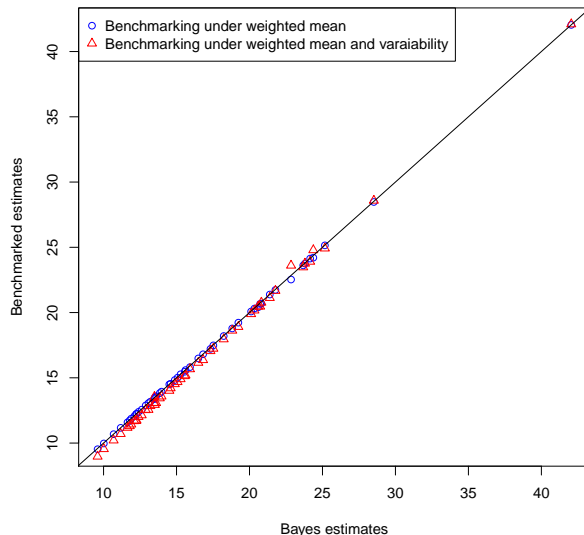


Figure 1: Benchmarking the mean alone leads to little change from the Bayes estimates; benchmarking both the mean and variability has very little improvement.

In each case, the benchmarking weights w_i are proportional to the estimated population of children aged 5–17 in each state. Intuitive remarks regarding how to choose h for benchmarking the weighted variability are given in [Datta et al. \(2011\)](#); [Ghosh and Steorts \(2013\)](#). We refer to these for technical details. Figure 1 compares the unconstrained Bayes estimator to the benchmarked Bayes estimators. Poverty estimates change very little when benchmarking the weighted mean alone, and only a little more when we benchmark both mean and variability.

The most important part of our procedure is picking the matrix Ω used to measure the smoothness of estimates—equivalently, picking the matrix Q which says how similar the estimates for any two domains should be. This is inevitably application-specific. In the results reported below, we used the simple choice where $q_{ii'} = 1$ if the states i and i' shared a border, and 0 otherwise. This treats the states as nodes in an unweighted graph, with Q being its adjacency matrix and Ω its Laplacian. As described in §3.3, the smoothing factor γ was picked by leave-one-out cross-validation; the final

value was $\gamma \approx 0.02$. Figure 4 shows the smoothed and mean-benchmarked Bayes estimates versus the unconstrained Bayes estimates. In general, low Bayes estimates are pulled up and high ones are brought down; everywhere, estimates are adjusted towards their neighbors.

Figure 2 further illustrates the effects of combining smoothing with benchmarking a weighted mean. It is broken up by the four U.S. Census regions⁶ for ease of visualization. While benchmarking alone has relatively little impact on the Bayes estimates, benchmarking plus smoothing does. In each region, the smoothed estimates fall on lines of slope less than 1, indicating shrinkage towards a common value, even though the regions are not part of the smoothing scheme. This means that the value toward which the estimates are shrunk is not necessarily the regional mean—observe region 2, where most constrained estimates *exceed* the Bayes estimates.

Figure 3 shows the statewise MSEs under the bootstrap of §4 for different combinations of benchmarking and smoothing. Smoothing tends to bring down the bias and the MSE for most but not all states—it is in fact known that the bias cannot be reduced uniformly across areas (Ghosh and Steorts 2013; Pfeiffermann 2013). The one “state” for which smoothing drastically increases the estimated error is the District of Columbia, which is unsurprising on substantive grounds.⁷

We considered several alternative ways of smoothing the Bayes estimates. One was to make $q_{ii'}$ decrease with the geographic distance between states, regarded as points at either their centers or their capitals. However, neither choice of representative point was compelling, and we would also have to pick the exact rate at which $q_{ii'}$ decreased with distance. A second approach was to treat the Census regions as clusters, setting $q_{ii'} = 1$ within a cluster and $q_{ii'} = 0$ between them. This however performed poorly un-

⁶Region 1, the Northeast: Connecticut, Maine, Massachusetts, New Hampshire, New Jersey, New York, Pennsylvania, Rhode Island, Vermont. Region 2, the Midwest: Illinois, Indiana, Iowa, Kansas, Michigan, Minnesota, Missouri, Nebraska, North Dakota, Ohio, South Dakota, Wisconsin. Region 3, the South: Alabama, Arkansas, Delaware, the District of Columbia, Florida, Georgia, Kentucky, Louisiana, Maryland, Minnesota, North Carolina, Oklahoma, South Carolina, Tennessee, Texas, Virginia, West Virginia. Region 4, the West: Arizona, California, Colorado, Idaho, Montana, Nevada, New Mexico, Oregon, Utah, Washington, Wyoming. Note that Alaska and Hawaii are not included in this regional scheme.

⁷Briefly, the District of Columbia is not a separate state, but rather a small part of a larger metropolitan area, containing a disproportionate share of the metropolis’s poorest neighborhoods. The adjoining states of Maryland and Virginia are much larger, much more prosperous, and much more heterogeneous. Many of these issues would be alleviated if we had data on finer spatial scales.

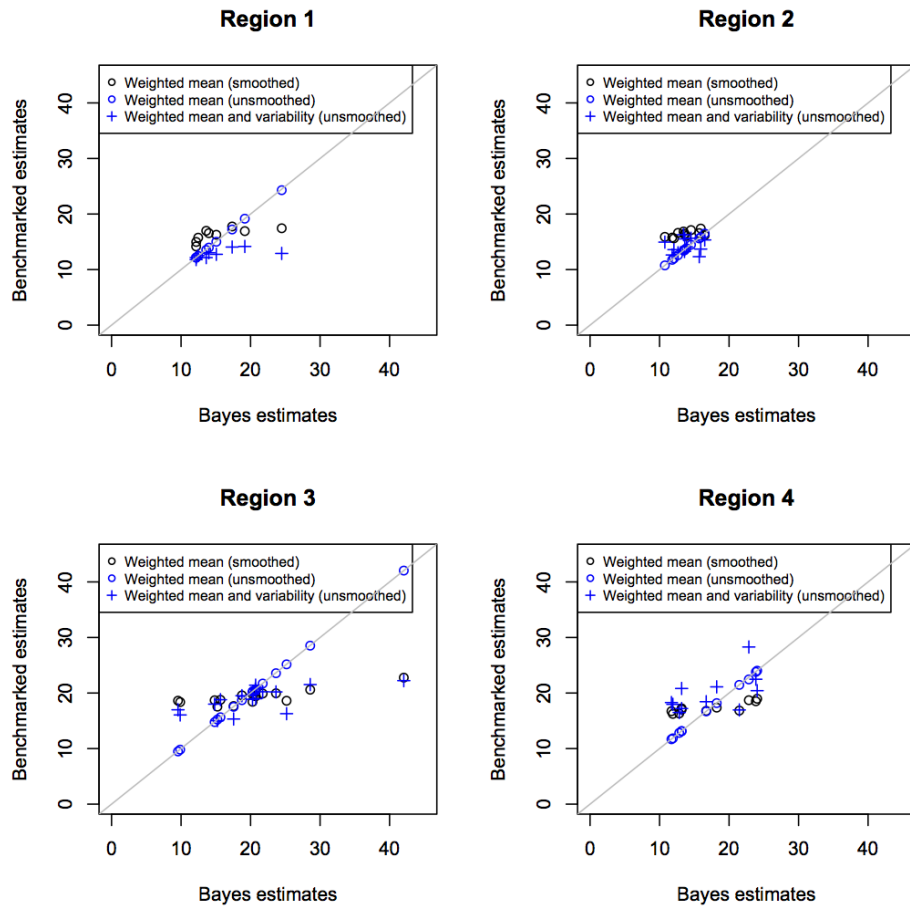


Figure 2: Benchmarked estimates with and without region, plotted against the Bayes estimates, by region.

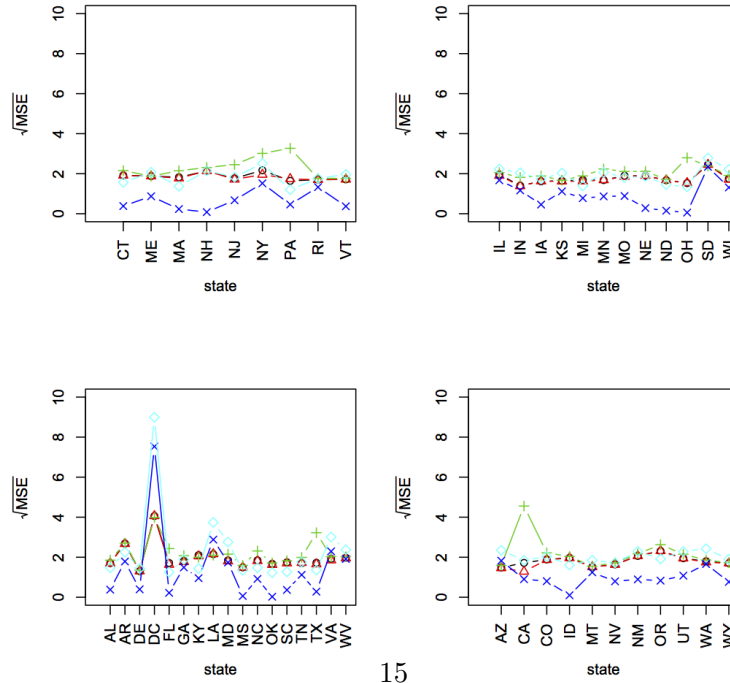
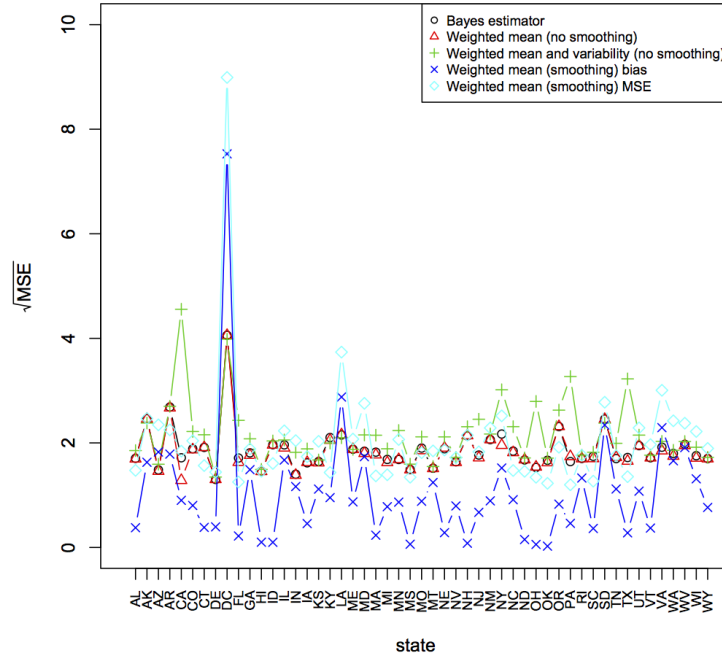


Figure 3: Above: Bootstrap MSEs for the SAIPE data and the Fay-Harriot model, under different combinations of benchmarking and smoothing. Below: the same data, but broken into the geographic regions.

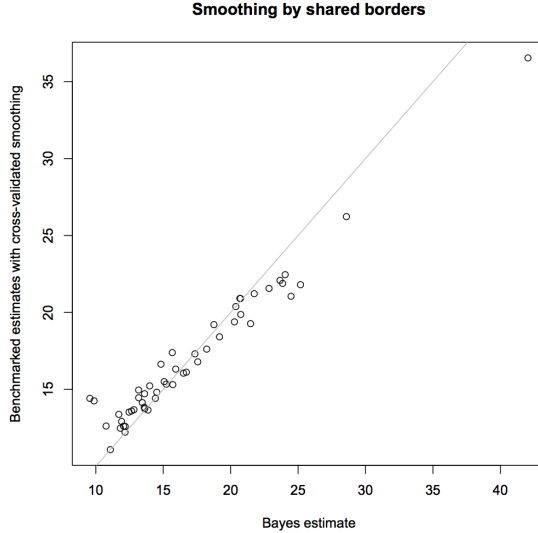


Figure 4: Smoothed, mean-constrained Bayes estimates versus unconstrained Bayes estimates.

der cross-validation. Similarly, we attempted diffusion-map k -means⁸ (Lee and Wasserman 2010; Richards 2014) with varying k , but none worked well under cross-validation.

We suspect that both distance-based and clustering approaches may work better at finer levels of spatial resolution, e.g., moving from whole states to counties or even census tracts, which are more demographically homogeneous. However, this is speculative without such fine-grained data.

6 Discussion

We have provided a general approach to SAE at both the unit and area levels, where we smooth and benchmark estimates. Our approach yields closed-form solutions without requiring any distributional assumptions. Furthermore, our results apply for linear and non-linear estimators and extend to multivariate settings. Finally, we show through a bootstrap approxima-

⁸The covariates were the state mean income, the state median income, the fraction of adults with at least a high-school education, the percentage of the population racially classified as white, and the percentage living in metropolitan areas.

tion and cross-validation that smoothing does improve estimation of poverty rates at the state level for the SAIPE dataset for most states as measured by MSE.

Note that we do not provide a simulation study, since any one that we considered was an unfair and biased comparison to either our proposed estimator or those proposed earlier in the literature. This is due to the fact that the Fay-Herriot model does not assume a smoothing/spatial component, however, our loss function does. We would be glad to consider any fair simulation study if one can be pointed out that would lead to helpful and meaningful results.

Another direction for future research is the extension of the present work to address weighted variability constraints as well. This becomes a difficult non-convex optimization problem, for which it is not clear how to efficiently and reliably obtain a numerical solution. The imposition of more than one weighted variability constraint specifies the feasible set as the intersection of multiple $(m - 1)$ -dimensional manifolds in m -dimensional Euclidean space. Careful consideration of the geometry of the resulting optimization problem may yield insight into methods of obtaining exact or approximate solutions, at least in certain special cases. Such ideas are clearly a potential direction for future work.

Throughout, we have worked with squared error. However, it should be possible to replace this with any other convex norm, with minimal changes to our approach. Once the Bayes estimate is obtained, the constrained Bayes estimate would be found by projection onto the feasible set. This would presumably mean more numerical optimization and fewer closed forms, but the optimization would remain convex and tractable. Getting the initial Bayes estimates under a different loss function might be more challenging.

It may be possible to go beyond point estimates to distributional estimates. Given a sample from the posterior distribution (e.g., from MCMC), it is possible to project each sample point into the feasible set, giving a posterior distribution whose support respects the constraints. The inferential validity of this sample would however require careful investigation.⁹

Acknowledgements

RCS was supported by NSF grants SES1130706 and DMS1043903 and NIH grant #1 U24 GM110707-01.

⁹Note that this is rather different from the idea in the recent papers of [Zhu et al. \(2012\)](#) of regularizing the posterior distribution, where the constraints or penalties are expressed as functionals of the whole posterior distribution.

References

- Battese, G., Harter, R., and Fuller, W. (1988), “An Error-Components Model for Prediction of County Crop Area Using Survey and Satellite Data,” *Journal of the American Statistical Association*, 83, 28–36.
- Belkin, M., Niyogi, P., and Sindhvani, V. (2006), “Manifold Regularization: A Geometric Framework for Learning from Labeled and Unlabeled Examples,” *Journal of Machine Learning Research*, 7, 239–2434.
- Bell, W., Datta, G., and Ghosh, M. (2013), “Benchmarked Small Area Estimators,” *Biometrika*, 100, 189–202.
- Corona, E., Lane, T., Storlie, C., and Neil, J. (2008), “Using Laplacian Methods, RKHS Smoothing Splines and Bayesian Estimation as a framework for Regression on Graph and Graph Related Domains,” Tech. Rep. TR-CS-2008-06, Department of Computer Science, University of New Mexico.
- Datta, G. and Ghosh, M. (1991), “Bayesian Prediction in Linear Models: Applications to Small Area Estimation,” *The Annals of Statistics*, 19, 1748–1770.
- Datta, G. S., Ghosh, M., Steorts, R., and Maples, J. (2011), “Bayesian benchmarking with applications to small area estimation,” *TEST*, 20, 574–588.
- Davison, A. C. and Hinkley, D. V. (1997), *Bootstrap Methods and their Applications*, Cambridge, England: Cambridge University Press.
- Fay, R. and Herriot, R. (1979), “Estimates of income from small places: an application of James-Stein procedures to census data,” *Journal of the American Statistical Association*, 74, 269–277.
- Ghosh, M. (1992), “Constrained Bayes estimation with applications,” *Journal of the American Statistical Association*, 87, 533–540.
- Ghosh, M. and Rao, J. (1994), “Small area estimation: an appraisal,” *Statistical science*, 55–76.
- Ghosh, M. and Steorts, R. C. (2013), “Two-Stage Bayesian Benchmarking as Applied to Small Area Estimation,” *TEST*, 22, 670–687.

- Kafadar, K. (1996), “Smoothing Geographical Data, Particularly Rates of Disease,” *Statistics in Medicine*, 15, 2539–2560.
- Lee, A. B. and Wasserman, L. (2010), “Spectral Connectivity Analysis,” *Journal of the American Statistical Association*, 105, 1241–1255.
- Louis, T. (1984), “Estimating a population of parameter values using Bayes and empirical Bayes methods,” *Journal of the American Statistical Association*, 79.
- Newman, M. E. J. (2010), *Networks: An Introduction*, Oxford, England: Oxford University Press.
- Pfeffermann, D. (2013), “New important developments in small area estimation,” *Statistical Science*, 28, 40–68.
- Pratesi, M. and Salvati, N. (2008), “Small area estimation: the EBLUP estimator based on spatially correlated random area effects,” *Statistical methods and applications*, 17, 113–141.
- Rao, J. (2003), *Small Area Estimation*, Wiley, New York.
- Richards, J. (2014), *diffusionMap*, r package version 1.1-0.
- Souza, D. F., Moura, F., and Migon, H. (2009), “Small area population prediction via hierarchical models,” *Catalogue no. 12-001-X*, 203.
- Stone, M. (1974), “Cross-validatory choice and assessment of statistical predictions,” *Journal of the Royal Statistical Society B*, 36, 111–147.
- Wahba, G. (1990), *Spline Models for Observational Data*, Philadelphia: Society for Industrial and Applied Mathematics.
- Wehbe, L., Ramdas, A., Steorts, R. C., and Shalizi, C. R. (2014), “Regularized Brain Reading with Shrinkage and Smoothing,” *Annals of Applied Statistics*, submitted.
- Zhu, J., Chen, N., and Xing, E. P. (2012), “Bayesian Inference with Posterior Regularization and Infinite Latent SVMs,” *arXiv preprint arXiv:1210.1766*.

A Results for Unit-Level Models

Many problems feature multiple levels of aggregation. For simplicity, we consider the specific case of two levels (from which the extension to three or more levels will be fairly clear). “Areas” refer to the upper level of aggregation and are divided into **units**. The i th area contains n_i units; the total number of units is $N = \sum_i n_i$. Units are strictly nested within areas and are indexed by j . We denote the area-level quantities as θ_i^A (with covariates \mathbf{x}_i^A , etc.), and the unit-level parameters as θ_{ij}^U (with covariates \mathbf{x}_{ij}^U , etc.). Denote the vectors of Bayes estimates by $\hat{\boldsymbol{\theta}}_A^B$ and $\hat{\boldsymbol{\theta}}_U^B$. The loss weight for unit j in area i is ξ_{ij} . Assume that loss is additive across areas and units; thus, the total loss from the action (estimate) $(\boldsymbol{\delta}^A, \boldsymbol{\delta}^U)$ is

$$\sum_i \phi_i (\delta_i^A - \theta_i^A)^2 + \sum_{ij} \xi_{ij} (\delta_{ij}^U - \theta_{ij}^U)^2.$$

Define Ξ as the diagonal matrix of the ξ_{ij} , which is positive-definite.

In many important cases, the area-level parameters are functions (e.g., weighted means or proportions) of the parameters for the units contained within the area (e.g., we might use $\bar{\theta}_{iw} = \sum_j w_{ij} \theta_{ij}^U$ as our θ_i^A). Less trivial examples are quantiles or Gini coefficients of the θ_{ij}^U for each area. However, it does not make sense for the unit-level parameters to be functions of the area-level parameters. The area-level parameter does not *have* to be a function of the unit-level parameters (e.g., if we have random effects for both areas and for units, the latter do not determine the former).

The general results of Theorems 1 and 2 can also be applied to models at the unit level, as described below.

A.1 Smoothing for Unit-Level Models

Consider the case where each area is partitioned into units, and estimates are sought at both the unit and the area level. (See §1.1 for notation.) We need two similarity functions, q_A as before, and q_U , where $q_U(x_{ij}, x_{i'j'})$ is the similarity between unit j in area i and unit j' in area i' . The smoothness-

augmented loss function is

$$\begin{aligned}
L(\boldsymbol{\theta}^A, \boldsymbol{\theta}^U, \boldsymbol{\delta}^A, \boldsymbol{\delta}^U) &= \sum_i \phi_i (\delta_i^A - \theta_i^A)^2 + \sum_{ij} \xi_{ij} (\delta_{ij}^U - \theta_{ij}^U)^2 \\
&\quad + \gamma_A \sum_{i,i'} (\delta_i^A - \delta_{i'}^A)^2 q_{ii'}^A + \gamma_U \sum_{ij,i'j'} (\delta_{ij}^U - \delta_{i'j'}^U)^2 q_{ij,i'j'}^U \\
&= (\boldsymbol{\delta}_A - \boldsymbol{\theta}_A)^T \Phi (\boldsymbol{\delta}_A - \boldsymbol{\theta}_A) + (\boldsymbol{\delta}_U - \boldsymbol{\theta}_U)^T \Xi (\boldsymbol{\delta}_U - \boldsymbol{\theta}_U) \\
&\quad + \gamma_A \boldsymbol{\delta}_A \Omega_A \boldsymbol{\delta}_A + \gamma_U \boldsymbol{\delta}_U \Omega_U \boldsymbol{\delta}_U,
\end{aligned} \tag{8}$$

defining Ω_A and Ω_U via Lemma 1.

Corollary 1. *The posterior risk of the loss (8) is minimized by the estimators $\tilde{\boldsymbol{\theta}}_A^S = (I_m + \gamma_A \Phi^{-1} \Omega_A)^{-1} \hat{\boldsymbol{\theta}}_A^B$ and $\tilde{\boldsymbol{\theta}}_U^S = (I_N + \gamma_U \Xi^{-1} \Omega_U)^{-1} \hat{\boldsymbol{\theta}}_U^B$.*

Proof. First, note that the “ m ” of Theorem 1 is in fact $m + N$ in this setting. Now partition $\boldsymbol{\theta} = (\boldsymbol{\theta}^A, \boldsymbol{\theta}^U)$. Similarly, partition the estimate vector as $\tilde{\boldsymbol{\theta}}^S = (\tilde{\boldsymbol{\theta}}_A^S, \tilde{\boldsymbol{\theta}}_U^S)$. Set both the Φ and Ω matrices to be block-diagonal:

$$\Phi = \begin{bmatrix} \Phi & 0 \\ 0 & \Xi \end{bmatrix}, \quad \Omega = \begin{bmatrix} \Omega_A & 0 \\ 0 & \frac{\gamma_U}{\gamma_A} \Omega_U \end{bmatrix}.$$

Now apply Theorem 1. □

Remark. This device of partitioning was employed by Datta and Ghosh (1991). It can be combined with multivariate parameters¹⁰, and indeed with more than two levels of spatial hierarchy, if needed. Since Φ and Ω are block-diagonal, the optimizations over $\tilde{\boldsymbol{\theta}}_A^S$ and $\tilde{\boldsymbol{\theta}}_U^S$ can be done separately, but no separate theorem is required.

A.2 Benchmarking for Unit-Level Models

Unit-level models can be benchmarked either for weighted means or for both weighted means and weighted variability.

A.2.1 Weighted Mean

Consider a unit-level model in which we wish to benchmark both the weighted mean of the area-level estimates and the weighted means of the unit-level

¹⁰As before, group the parameters in $\boldsymbol{\theta}^A$ and $\boldsymbol{\theta}^U$ by component. Then Φ and Ξ are diagonal; Ω_A and Ω_U are block-diagonal, each block a copy of the univariate Ω_A or Ω_U .

estimates within each area. Then we wish to minimize (8) under the constraints

$$\sum_i \eta_i \delta_i = t^A, \quad \sum_j w_{ij} \hat{\theta}_{ij}^U = \delta_i \quad \forall i. \quad (9)$$

Partition $\boldsymbol{\theta}$, $\hat{\boldsymbol{\theta}}^B$, and $\boldsymbol{\delta}$ as in §A.1. Define \tilde{W} as the $m \times N$ matrix such that¹¹ $(\tilde{W}\boldsymbol{\theta}^U)_i = \sum_j w_{ij} \theta_{ij}^U$, and define

$$M = \begin{bmatrix} \boldsymbol{\eta}^T & \mathbf{0}_N \\ -I_m & \tilde{W} \end{bmatrix},$$

where I_m is the $m \times m$ identity matrix and $\mathbf{0}_N$ is the length- N vector of zeroes. Let $\mathbf{t} = (t^A, \mathbf{0}_m)$, where again $\mathbf{0}_m$ is the length- m vector of zeroes. Then (9) amounts to $M\boldsymbol{\delta} = \mathbf{t}$. By a direct application of Theorem 2, we have the following result.

Corollary 2. *The benchmarked Bayes estimator that minimizes the posterior risk in (8) under the constraints in (9) is*

$$\tilde{\boldsymbol{\theta}}^S = \Sigma^{-1} \left[\Phi \hat{\boldsymbol{\theta}}^B + M^T (M \Sigma^{-1} M^T)^{-1} (\mathbf{t} - M \Sigma^{-1} \Phi \hat{\boldsymbol{\theta}}^B) \right],$$

where $\Sigma = \Phi + \gamma \Omega$, and where Φ and Ω are as in the proof of Corollary 1.

B Lemma on Squared Differences

Lemma 1. *For a suitable matrix Ω ,*

$$\sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'} = \boldsymbol{\delta}^T \Omega \boldsymbol{\delta}.$$

Proof. Begin by expanding the square and collecting terms:

$$\begin{aligned} & \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'} \\ &= \sum_{i,i'} \delta_i^2 q_{ii'} + \sum_{i,i'} \delta_{i'}^2 q_{ii'} - 2 \sum_{i,i'} \delta_i \delta_{i'} q_{ii'} \\ &= \sum_i \delta_i^2 \sum_{i'} q_{ii'} + \sum_{i'} \delta_{i'}^2 \sum_i q_{ii'} - 2 \sum_{i,i'} \delta_i \delta_{i'} q_{ii'}. \end{aligned}$$

¹¹The i^{th} row of \tilde{W} will have non-zero entries w_{ij} in the columns corresponding to the units in area i , and zeroes everywhere else.

Now define the diagonal matrix $Q^{(r)}$ with elements $q_{ii}^{(r)} = \sum_{i'} q_{ii'}$, and define the diagonal matrix $Q^{(c)}$ with elements $q_{jj}^{(c)} = \sum_i q_{ij}$. Substituting,

$$\begin{aligned} \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{i,i'} &= \boldsymbol{\delta}^T Q^{(r)} \boldsymbol{\delta} + \boldsymbol{\delta}^T Q^{(c)} \boldsymbol{\delta} - 2\boldsymbol{\delta}^T Q \boldsymbol{\delta} \\ &= \boldsymbol{\delta}^T \left(Q^{(r)} + Q^{(c)} - 2Q \right) \boldsymbol{\delta}, \end{aligned}$$

which defines Ω . □

Remark. In an unweighted, undirected graph with adjacency matrix A , the degree matrix D is defined by $D_{ii} = \sum_j A_{ij}$, $D_{ij} = 0$; the graph Laplacian in turn is $L = D - A$ (Newman 2010). If Q is an adjacency matrix, then $Q^{(r)} = Q^{(c)} = D$, and $\Omega = 2L$.

Remark. By construction, Ω is clearly positive semi-definite. It is not positive definite, because $(1 \ 1 \ \dots \ 1)$ is always an eigenvector, of eigenvalue zero. This corresponds to the fact that adding the same constant to each δ_i does not change $\sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{i,i'}$. (These are of course basic properties of graph Laplacians.)

C Residual Bootstrap

We consider the model

$$\begin{aligned} y_i &= \theta_i + U_i \\ \theta_i &= \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i \end{aligned}$$

where $i = 1, \dots, m$ and where the observational noise vector \mathbf{U} has a known diagonal covariance matrix Σ_U , with the i th diagonal element of Σ_U denoted by $\sigma_{U,i}^2$. We impose constraints (benchmarking, smoothing) on the estimates of the θ_i in order to better regularize and borrow strength. Call the constrained estimates $\tilde{\boldsymbol{\theta}}^{BM}$. Then we can define residuals for each observation:

$$r_i = y_i - \tilde{\theta}_i^{BM}.$$

If we standardize these as

$$\tilde{r}_i = \frac{y_i - \tilde{\theta}_i^{BM}}{\sigma_{U,i}},$$

we get quantities which should have the same distribution for all areas, if our constraints are valid and our model fits well. We then bootstrap by re-sampling these residuals:

$$u_i^* \stackrel{\text{iid}}{\sim} \tilde{\mathbf{r}}$$

$$y_i^* = \theta_i^{BM} + u_i^* \sigma_{U,i}$$

where $i = 1, \dots, m$. Note that the first line of the above model simply means that we draw iid random variables u_1^*, \dots, u_m^* where each u_i^* is equal to each of $\tilde{r}_1, \dots, \tilde{r}_m$ with probability $1/m$. Re-sampling-based bootstraps are commonly used in assessing uncertainty for regression models. They presume the correctness of the functional form of the regression, but not of distributional assumptions about the noise.¹²

To summarize, the resampling procedure would be this:

1. From data (\mathbf{x}, \mathbf{y}) , obtain constrained estimates $\tilde{\boldsymbol{\theta}}^{BM}$ and residuals $\mathbf{r} = \mathbf{y} - \tilde{\boldsymbol{\theta}}^{BM}$.
2. Calculate standardized residuals $\tilde{\mathbf{r}} = \Sigma_U^{-1/2} \mathbf{r}$.
3. Repeat B times:
 - (a) Draw \mathbf{u}^* by resampling with replacement from $\tilde{\mathbf{r}}$.
 - (b) Set $\mathbf{y}^* = \tilde{\boldsymbol{\theta}}^{BM} + \Sigma_U^{-1/2} \mathbf{u}^*$.
 - (c) Re-run inference on $(\mathbf{x}, \mathbf{y}^*)$ to get $\tilde{\boldsymbol{\theta}}^{BM*}$.
4. Use the distribution of $\tilde{\boldsymbol{\theta}}^{BM*}$ in bootstrap calculations.

¹²There is also a “wild bootstrap” (Davison and Hinkley 1997, p. 272) which would evade having to know the observational noise variances, at some cost in efficiency.